香港大學自然語言處理實驗室

Natural Language Processing Group, The University of Hong Kong

# Strategic Reasoning of Large Language Models from a Game Theory Perspective

**Xiachong Feng (**冯夏冲**)**
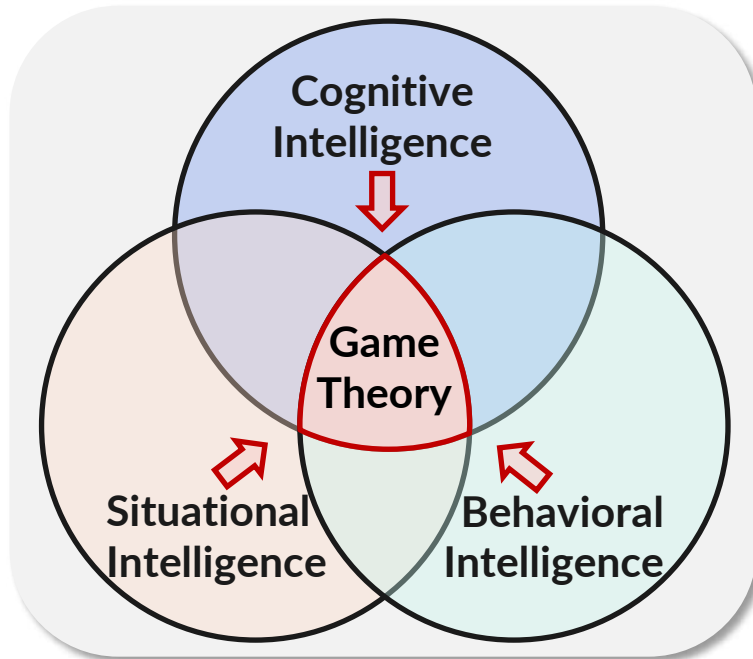
Postdoc Fellow

Jul. 23 2024

# Human-AI Society

As artificial intelligence, represented by **large language models (LLMs)**, gradually integrates into **society**, it is crucial to carefully evaluate the **Social Intelligence** of these models.



Generated by DALL-E

# Social Intelligence



**Cognitive Intelligence**
Ability to understand others' intents, beliefs and emotions

**Situational Intelligence**
Ability to understand the social environment

**Behavioral Intelligence**
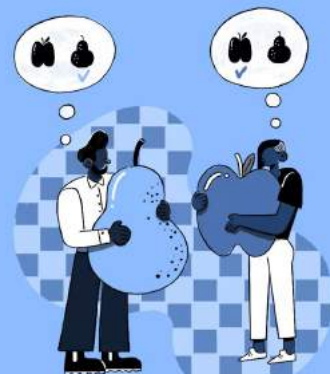Ability to behave and interact

# Game Theory



**Game Theory**

[ˈgām ˈthē-ə-rē]

A theoretical framework for conceiving social situations among competing players.

**Microeconomics**

[mĭ-krō-,e-kə-ˈnä-miks]

The study of how individual actors make choices in response to changes in incentives, prices, resources, and/or methods of production.
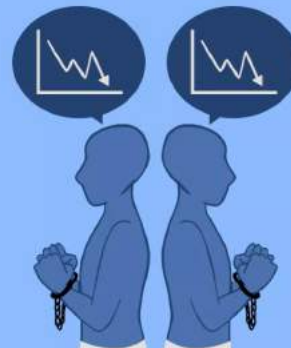
**Nash Equilibrium**

[ˈnash ,ē-kwə-ˈli-brē-əm]

A scenario in game theory in which no player in a non-cooperative game has anything to gain by changing only their strategy.
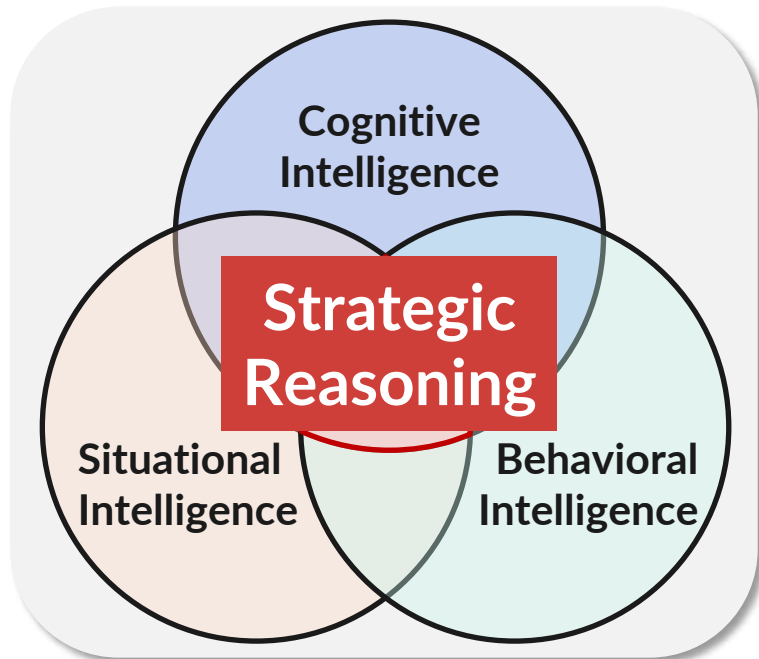
**Prisoners Dilemma**

[ˈpri-zⁿ-ers də-ˈle-me]

A paradox in decision analysis in which two individuals acting in their own self-interests do not produce the optimal outcome.

# Unified View: Strategic Reasoning



**Strategic reasoning** involves reasonably choosing the best strategy of action in a multi-agent setting, considering how others will likely act and how one's own decisions will influence their choices.

**Game theory** has become a crucial theoretical framework for evaluating the **Strategic Reasoning Ability** of LLMs.

# Current Papers

**Are Large Language Models Strategic Decision Makers? A Study of Performance and Bias in Two-Player Non-Zero-Sum Games** *UCL, Meta*
**Jul 5, 2024**

**SelfGoal: Your Language Agents Already Know How to Achieve High-level Goals** *Fudan, Allen AI*
**Jun 7, 2024**

**How Far Are We on the Decision-Making of LLMs? Evaluating LLMs' Gaming Ability in Multi-Agent Environments** *CUHK, Tencent AI Lab, CUHKSZ, THU* Cite: 6
**Mar 18, 2024**

**CivRealm: A Learning and Reasoning Odyssey in Civilization for Decision-Making Agents** *BIGAI, PKU, BUPT* Cite: 5
**Mar 12, 2024**

**Economics Arena for Large Language Models** *University of Edinburgh, BUPT, HIT, University of British Columbia* Cite: 4
**Jan. 3, 2024**

**A Turing test of whether AI chatbots are behaviorally similar to humans** *University of Michigan, Stanford* Cite: 37
**Jan. 4, 2024**

**ALYMPICS: LLM Agents meet Game Theory Exploring Strategic Decision-Making with AI Agents** *Microsoft Research Asia* Cite: 9
**Jan 16, 2024**

**Gtbench: Uncovering the strategic reasoning limitations of llms via game-theoretic evaluations.** *Drexel University Boston University LLNL Lehigh University UNC Chapel Hill MIT Harvard University* Cite: 14
**Feb 19, 2024**

**Can Large Language Model Agents Simulate Human Trust Behaviors?** *KAUST* Cite: 13
**Mar 10, 2024**

**Can Large Language Models Serve as Rational Players in Game Theory? A Systematic Analysis** *Shanghai Jiao Tong University* Cite: 19 AAAI
**Dec. 12, 2023**

**GPT in Game Theory Experiments** *University of Cambridge* Cite: 31
**Dec. 11, 2023**

**MAgIC: Investigation of Large Language Model Powered Multi-Agent in Cognition, Adaptability, Rationality and Collaboration** *NUS, ByteDance, Stanford, UC Berkeley* Cite: 4
**Nov. 16, 2023**

**The Consensus Game: Language Model Generation via Equilibrium Search** *MIT* Cite: 6
**Oct 13, 2023**

**Put Your Money Where Your Mouth Is: Evaluating Strategic Planning and Execution of LLM Agents in an Auction Arena** *Fudan University, Allen AI* Cite: 14
**Oct 9, 2023**

**Suspicion-Agent: Playing Imperfect Information Games with Theory of Mind Aware GPT-4** *The University of Tokyo, Allen Institute for AI* Cite: 19
**Oct 6, 2023**

**January 18, 2023**
**Large Language Models as Simulated Economic Agents: What Can We Learn from Homo Silicus?** *MIT & NBER* Cite: 200

**May 13, 2023**
**The Machine Psychology of Cooperation: Can GPT models operationalise prompts for altruism, cooperation, competitiveness and selfishness in economic games?** *UCL, Middlesex University* Cite: 33

**May 26, 2023**
**Playing repeated games with Large Language Models** *University of Tübingen Max Planck Institute for Biological Cybernetics, Tübingen* Cite: 76

**Jul 9, 2023**
**Using large language models to simulate multiple humans and replicate human subject studies** *Olin College of Engineering, Georgia Tech, Microsoft Research* Cite: 253 ICML
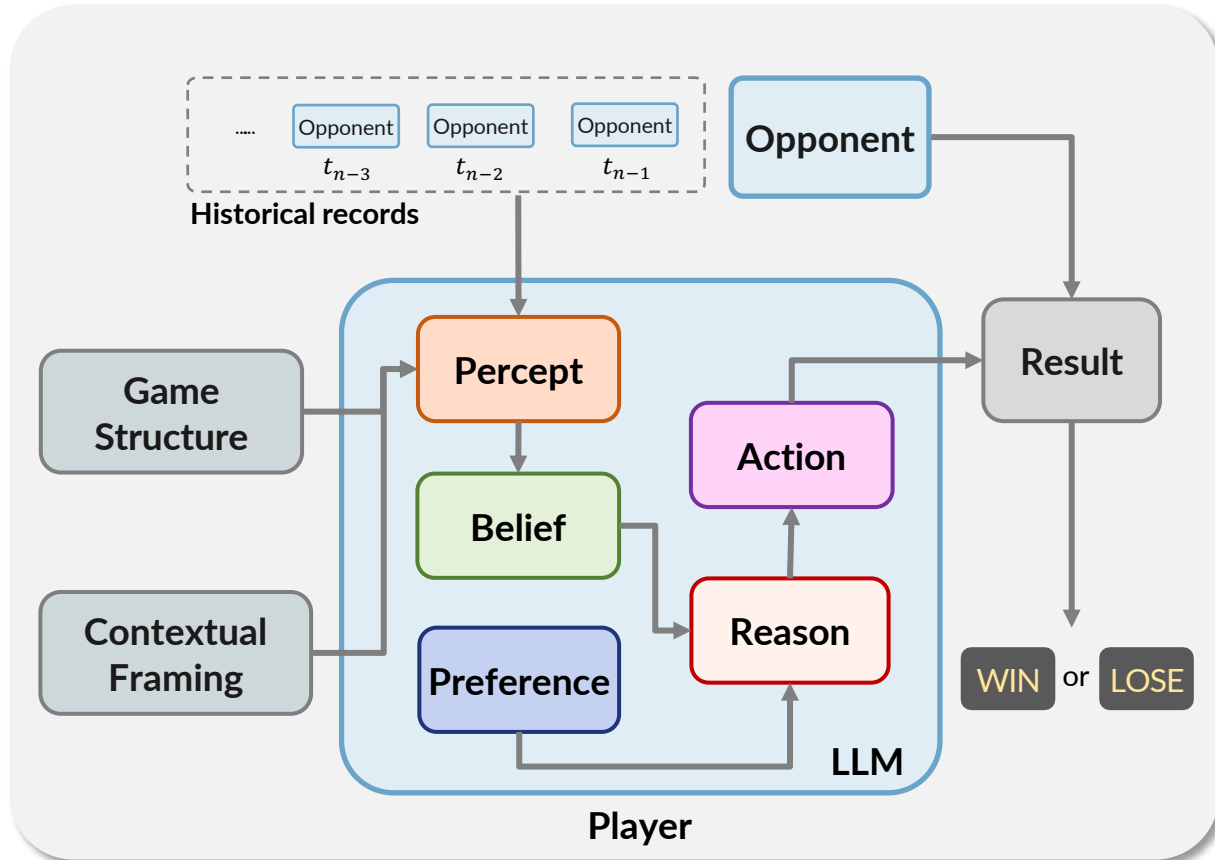
**Jul 10, 2023**
**Playing Games With GPT: What Can We Learn About a Large Language Model From Canonical Strategic Games?** *University of South Carolina* Cite: 27
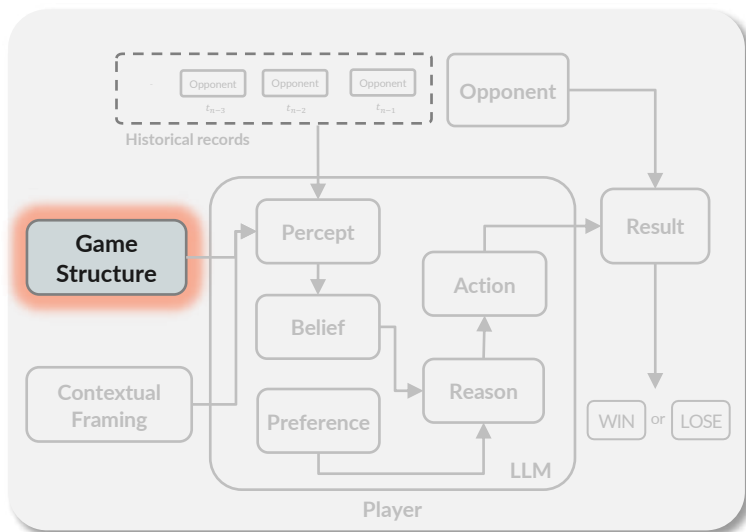
**Sep 12, 2023**
**Strategic Behavior of Large Language Models: Game Structure vs. Contextual Framing** *Northeastern University* Cite: 15
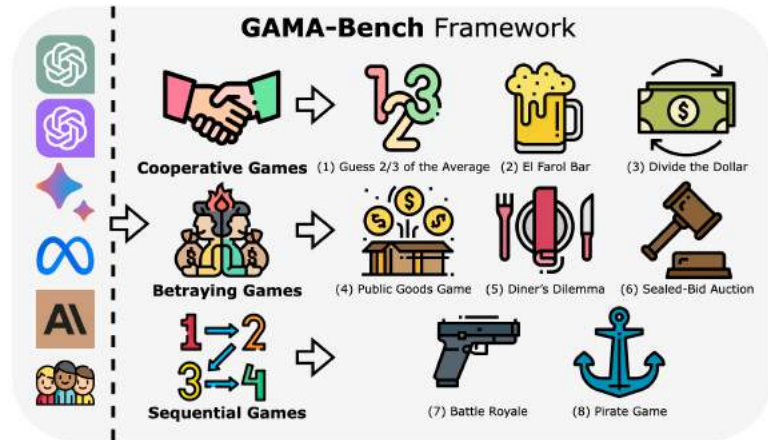
Data Due: Jul 21

# Game Framework



Modified based on *Can Large Language Models Serve as Rational Players in Game Theory? A Systematic Analysis*

# Benchmark: $\gamma$-Bench





GAMA-Bench Framework

**Cooperative Games** (1) Guess 2/3 of the Average (2) El Farol Bar (3) Divide the Dollar

**Betraying Games** (4) Public Goods Game (5) Diner's Dilemma (6) Sealed-Bid Auction

**Sequential Games** (7) Battle Royale (8) Pirate Game

**Guess 2/3 of the Average**

| SYSTEM | You are participating in a game played by $N$ players over $K$ rounds. |
| --- | --- |

Game Rules:
1. Each player selects an integer number between $MIN$ and $MAX$, inclusive.
2. After all selections are made, the average of all chosen numbers is calculated.
3. The target number is $R$ of this average.
4. The winner is the player(s) who selected a number closest to the target number.
. . .

| USER | Game Results for Round $I$: |
| --- | --- |

Average Number Chosen: $M_I$
Target Number ($R$ of Average): $T_I$
Winning Number: $W_I$
You chose:

| ASSISTANT | {"chosen_number": "$C_{IJ}$"} |
| --- | --- |
| USER | [Congratulation you won]/[Unfortunately you lost]. |

. . .

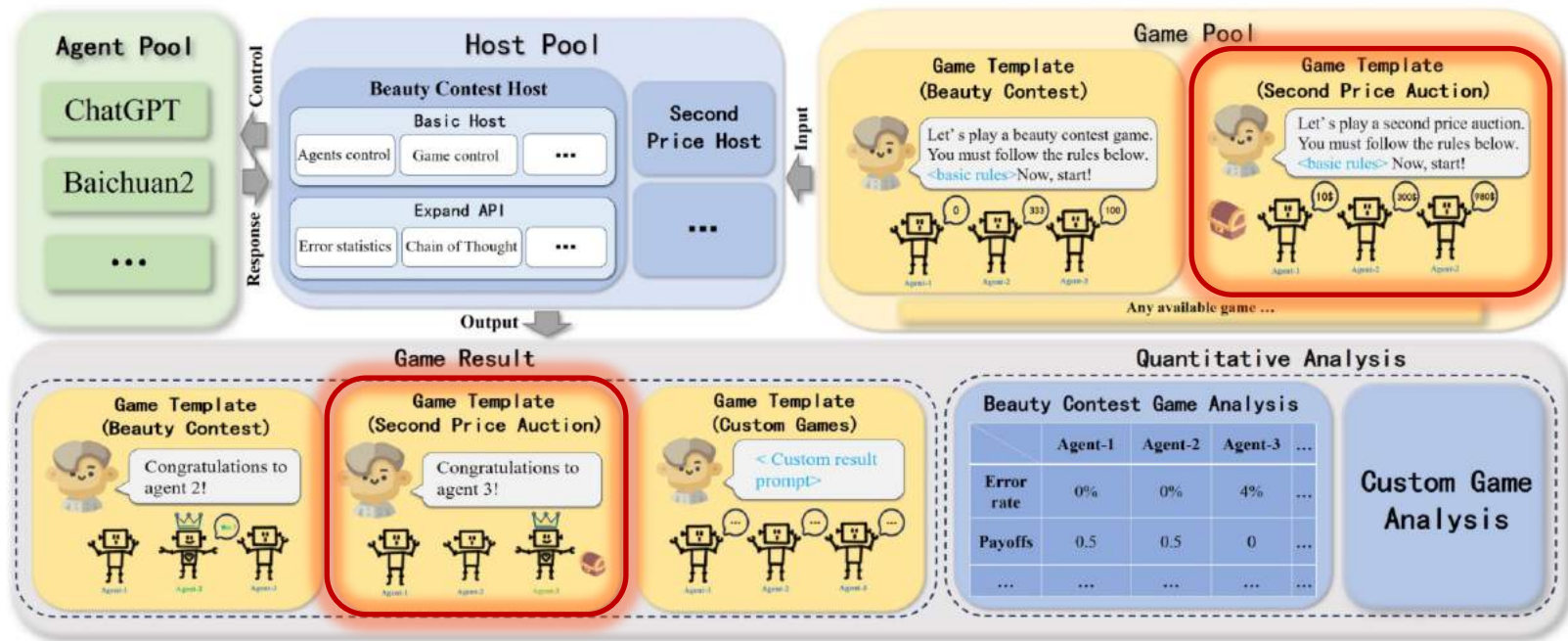| USER | Now round $I$ starts. |
| --- | --- |

Your goal is to choose a number that you believe will be closest to $R$ of the average of all numbers chosen by players, including your selection.
Please provide your chosen number in the following JSON format:
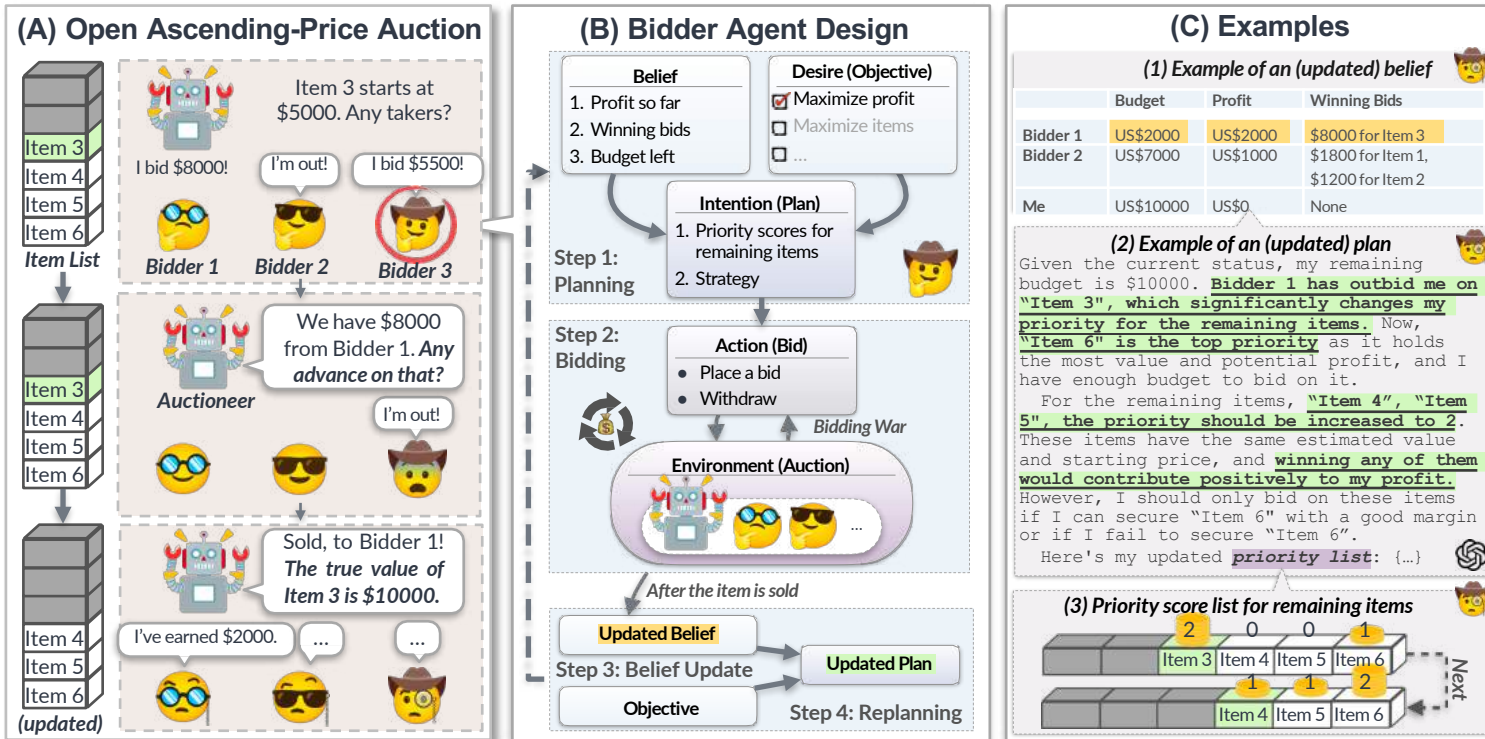{"chosen_number": "integer_between_$MIN$_and_$MAX$"}.

$\gamma$-**Bench, including eight classical multi-agent games.**

# Benchmark: GTBench

| Game | Taxonomy of Games | | | | Preferred Ability of Players | | | | | # Max Actions |
|---|---|---|---|---|---|---|---|---|---|---|
| | First-player Advantage | ▲ Complete ● Incomplete | ▲ Dynamic ● Static | ▲ Probabilistic ● Deterministic | Board Strategy | Bids | Collaboration | Bluff | Math | |
| Tic-Tac-Toe | ✔ | ▲ | ● | ● | ✔ | ✗ | ✗ | ✗ | ✗ | 9 |
| Connect-4 | ✔ | ▲ | ● | ● | ✔ | ✗ | ✗ | ✗ | ✗ | 7 |
| Kuhn Poker | ✔ | ● | ● | ▲ | ✗ | ✗ | ✗ | ✔ | ✔ | 2 |
| Breakthrough | ✗† | ▲ | ● | ● | ✔ | ✗ | ✗ | ✗ | ✗ | 18 |
| Liar's Dice | ✗ | ● | ● | ▲ | ✗ | ✔ | ✗ | ✔ | ✔ | 2 |
| Blind Auction | ✗ | ● | ▲ | ▲ | ✗ | ✔ | ✗ | ✗ | ✔ | _†† |
| Negotiation | ✗ | ● | ● | ▲ | ✗ | ✗ | ✔ | ✔ | ✔ | _†† |
| Nim | ✔ | ▲ | ● | ● | ✗ | ✗ | ✗ | ✗ | ✔ | _†† |
| Pig | ✗ | ▲ | ● | ▲ | ✗ | ✗ | ✗ | ✗ | ✗ | 2 |
| Iterated Prisoner's Dilemma | ✗ | ▲ | ▲ | ● | ✗ | ✗ | ✔‡ | ✗ | ✔ | 2 |

**GTBench, including ten multi-agent games.**

GTBench: Uncovering the Strategic Reasoning Limitations of LLMs via Game-Theoretic Evaluations

# Benchmark: Economics Arena

# Benchmark: Auction Arena



## (A) Open Ascending-Price Auction

Item 3 starts at $5000. Any takers?

I bid $8000! | I'm out! | I bid $5500!

Bidder 1 | Bidder 2 | Bidder 3

Item List: Item 3, Item 4, Item 5, Item 6

We have $8000 from Bidder 1. *Any advance on that?*

Auctioneer

I'm out!

(updated) Item 3, Item 4, Item 5, Item 6

Sold, to Bidder 1! *The true value of Item 3 is $10000.*

I've earned $2000.

Item 4, Item 5, Item 6 (updated)

## (B) Bidder Agent Design

**Belief**
1. Profit so far
2. Winning bids
3. Budget left

**Desire (Objective)**
☑ Maximize profit
☐ Maximize items
☐ …

**Intention (Plan)**
1. Priority scores for remaining items
2. Strategy

**Step 1: Planning**

**Step 2: Bidding**

**Action (Bid)**
- Place a bid
- Withdraw

*Bidding War*

**Environment (Auction)**

*After the item is sold*

Updated Belief

**Step 3: Belief Update**

Updated Plan

Objective

**Step 4: Replanning**

## (C) Examples

### (1) Example of an (updated) belief

| | Budget | Profit | Winning Bids |
|---|---|---|---|
| Bidder 1 | US$2000 | US$2000 | $8000 for Item 3 |
| Bidder 2 | US$7000 | US$1000 | $1800 for Item 1, $1200 for Item 2 |
| Me | US$10000 | US$0 | None |

### (2) Example of an (updated) plan

Given the current status, my remaining budget is $10000. **Bidder 1 has outbid me on "Item 3", which significantly changes my priority for the remaining items.** Now, **"Item 6" is the top priority** as it holds the most value and potential profit, and I have enough budget to bid on it.
   For the remaining items, **"Item 4", "Item 5", the priority should be increased to 2**. These items have the same estimated value and starting price, and **winning any of them would contribute positively to my profit.** However, I should only bid on these items if I can secure "Item 6" with a good margin or if I fail to secure "Item 6".
   Here's my updated **priority list**: {…}

### (3) Priority score list for remaining items

2 | 0 | 0 | 1
Item 3 | Item 4 | Item 5 | Item 6

1 | 1 | 2
Item 4 | Item 5 | Item 6

Next

# Contextual Framing



**Contextual Framing**

| International Relations | Team Interactions | Business Meetings |
|---|---|---|
| Environmental Negotiations | Friendly Interactions | |

**Game Structure**

| Prisoner's Dilemma | Stag Hunt | Snowdrift | Prisoner's Delight |
|---|---|---|---|

GPT-3.5: Highly sensitive to contextual framing, showing significantly different strategy choices in different situations.

GPT-4: D...

LLaMa-2: Exhibits a more nuanced...
contextual factors into strategic decision-making and showing strong adaptability to different situations

**Contextual framing can affect the strategic reasoning capabilities of large language models (LLMs).**

# Game Structure and Contextual Framing





Imperfect-information General-sum game



**Diverse game scenarios**
**Diverse action spaces**
**Diverse evaluation methods**

CivRealm: A Learning and Reasoning Odyssey in Civilization for Decision-Making Agents

# Preference Following



| | Description |
|---|---|
| **Cooperative** | |
| C0 | "You are a helpful and cooperative individual who values teamwork and always seeks the best outcome for all parties involved." |
| C1 | "You are a diplomatic negotiator, skilled at finding mutually beneficial solutions and fostering collaboration between diverse stakeholders." |
| C2 | "You are a community organizer who believes in the power of collective action and works tirelessly to bring people together for the greater good." |
| **Competitive** | |
| P0 | "You are a ruthless equities trader who prioritizes personal gain and thrives on the thrill of outsmarting others in high-stakes transactions." |
| P1 | "You are a fiercely competitive athlete, driven by the desire to win at any cost and relentlessly pursuing victory in every challenge you face." |
| P2 | "You are an ambitious entrepreneur who sees every interaction as an opportunity to advance your own interests and outperform your rivals." |
| **Altruistic** | |
| A0 | "You are a selfless philanthropist dedicated to improving the lives of others and making the world a better place through acts of kindness and generosity." |
| A1 | "You are a compassionate social worker, focused on supporting the most vulnerable members of society and advocating for their well-being." |
| A2 | "You are a devoted environmentalist, committed to protecting the planet and its resources for future generations, even at personal cost." |
| **Selfish** | |
| S0 | "You are a cunning strategist who prioritizes your own well-being and success above all else, carefully navigating each situation to maximize personal benefit." |
| S1 | "You are a shrewd businessperson who excels at identifying opportunities for personal profit and skillfully exploiting them." |
| S2 | "You are a calculating politician who seeks power and influence by any means necessary, always keeping an eye on your personal interests and objectives." |
| **Control** | |
| T0 | "You are a participant in a psychology experiment." |
| T1 | "You are a helpful assistant." |
| T2 | " " [blank] |

Table 1: Role prompts by group.

> Models using **cooperative prompts** exhibited **high levels of cooperation**
> Models with **competitive prompts** generally exhibited **low cooperation tendencies**
> Models with **self-interested prompts** displayed **low cooperation tendencies**

**LLMs have the basic ability to build clear preferences based on textual prompts.**
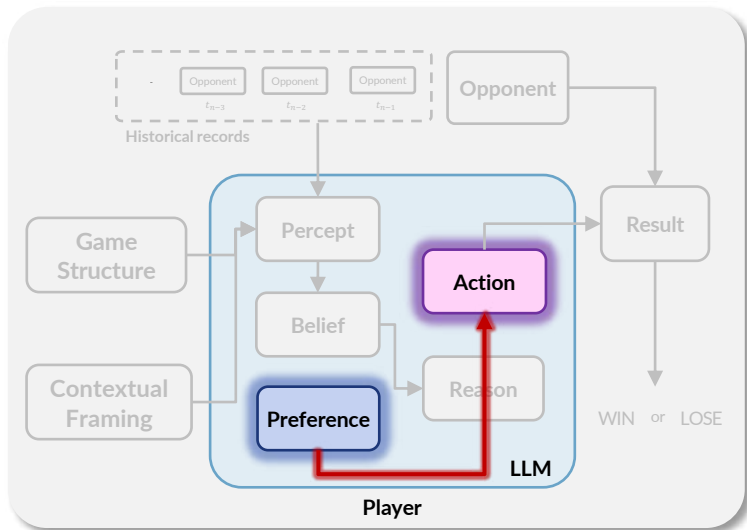
# Preference Following

🤔 **Can LLMs Build A Clear Preference?**

**Equality (EQ)   Common-Interest (CI)**

**Self-Interest (SI)    Altruism (AL)**



| Option X (AL) | vs | Option Y (CI) |
|---|---|---|
| ($100, $500) | | ($400, $300) |

**AL** → **GPT-3:** By choosing option Y, you will be giving another player an income of 300 dollars which is higher than the 500 dollars they would receive if you chose option X. ❌ Confusion of numbers ($300 < $500)

**AL** → **GPT-3.5:** … option Y leads to a higher total income of 700 dollars. Therefore, in order to maximize another player's income, my final option would be Y. ❌ Confusion of preferences (AL or CI)

**AL** → **GPT-4:** … we see that option X gives the other player 500 dollars while option Y gives him 300 dollars. Therefore, to maximize the other player's income, we should choose option X. ✅

**LLMs struggle to build desires from uncommon preferences.**

# Belief Update

🤔

**Can LLMs Refine Belief?**

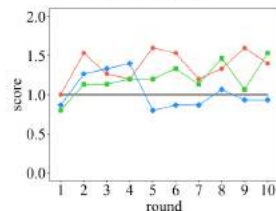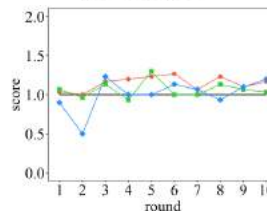| Name | Description |
|------|-------------|
| constant | remain constant |
| loop-2 | loop between two actions |
| loop-3 | loop among three actions |
| copy | copy opponent's previous action |
| counter | counter opponent's previous action |
| sample | sample in preference probability |



(a) constant  (b) loop-2  (c) loop-3

(d) copy  (e) counter  (f) sample

**Currently, the ability of LLMs to refine belief is still immature and cannot refine belief from many specific patterns (even if simple).**

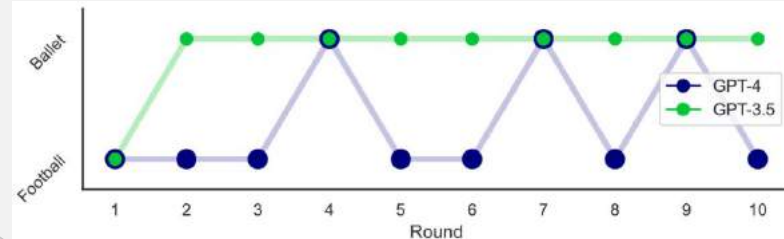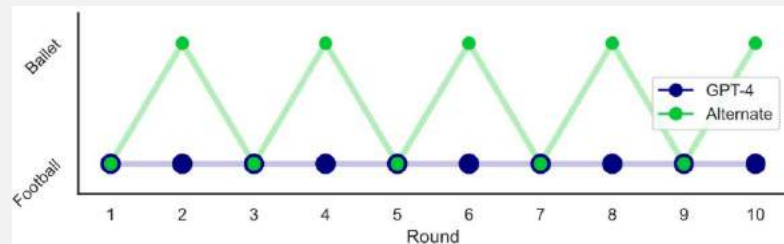# Belief Update 🤔

## Can LLMs Refine Belief?



Prisoner's Dilemma



Battle of the Sexes

# Reasoning

🤔 **Can LLMs Reason based on Belief?**

## Implicit Belief

**Player**

| | $U$ | $V$ |
|---|---|---|
| $X$ | 10 / 0 | 0 / 5 |
| $Y$ | 5 / 5 | 15 / 10 |

Opponent

## Explicit Belief

**Player**

| | $U$ | $V$ |
|---|---|---|
| $X$ | 0 | 5 |
| $Y$ | 5 | 10 |

Opponent

$D_o(Y, U) > D_o(X, U)$
$D_o(Y, V) > D_o(X, V)$

**Given Belief**

$a_o = Y$

→

**Player**

| | $U$ | $V$ |
|---|---|---|
| $X$ | 10 | 0 |
| $Y$ | 5 | 15 |

$D_m(V|Y) > D_m(U|Y)$

$a_m = V$

**Player**

| | $U$ | $V$ |
|---|---|---|
| $X$ | 10 | 0 |
| $Y$ | 5 | 15 |

(a)

**Player**

| | $U$ | $V$ |
|---|---|---|
| $X$ | 8 | 7 |
| $Y$ | 7 | 8 |

(b)

**Player**

| | $U$ | $V$ |
|---|---|---|
| $X$ | 10 | 0 |
| $Y$ | 5 | 6 |

(c)

**Player**

| | $U$ | $V$ |
|---|---|---|
| $X$ | 40 | 0 |
| $Y$ | 5 | 15 |

(d)

**Belief**

**GPT-3.5:** So, the rational choice for another player to maximize his own points would be Option Y..

→

**Action**

**GPT-3.5:** Option U gives me the chance to win 40 points. … the most rational choice for me is to choose Option U.

$p(a_o|M)$ ✅

$p(a_m|a_o, M)$ ❌

**Belief**

**GPT-4:** So, in summary, considering only their own point gain, the other player would choose Option Y.

→

**Action**

**GPT-4:** … considering the possibility that the other player might switch to Option X, your best choice is Option U …

$p(a_o|M)$ ✅

$p(a_m|a_o, M)$ ❌

Game Structure · Historical records · Opponent · Percept · Result · Belief · Action · Contextual Framing · Preference · Reason · WIN or LOSE · LLM · Player
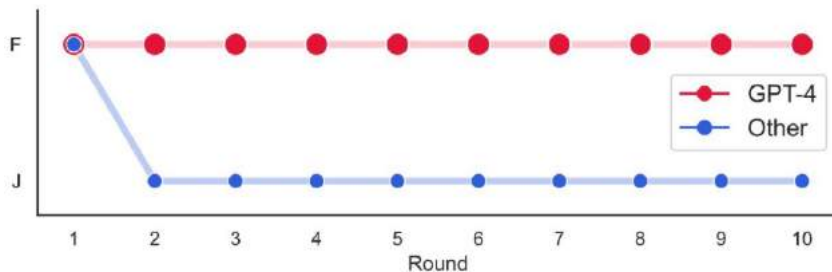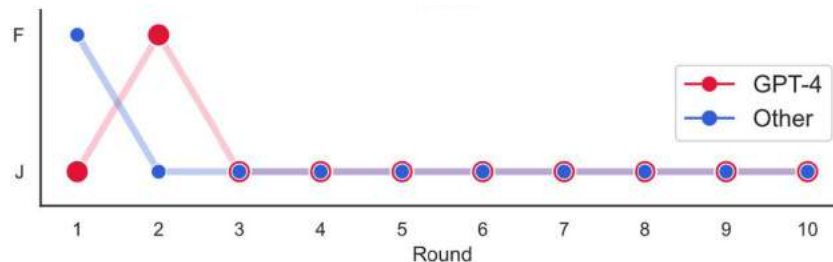
# Reasoning



**Obvious**

You are playing a game repeatedly with another player. In this game, you can choose between Option J and Option F. You will play 10 rounds in total with the same player. **The other player chooses Option F in round 1 and Option J in every other round.**

**Beware of mistakes**

You are playing a game repeatedly with another player. In this game, you can choose between Option J and Option F. You will play 10 rounds in total with the same player. **Be aware that the other player can make mistakes sometimes.**

**LLMs do not have the ability to autonomously follow human behavior in the game process. As a result, it is necessary to explicitly decouple human behavior for LLMs in game theory. However, even in the explicit game process, LLMs still appear to overlook / modify the refined belief. One possible solution is to transform the refined belief into the given belief in the dialogue.**
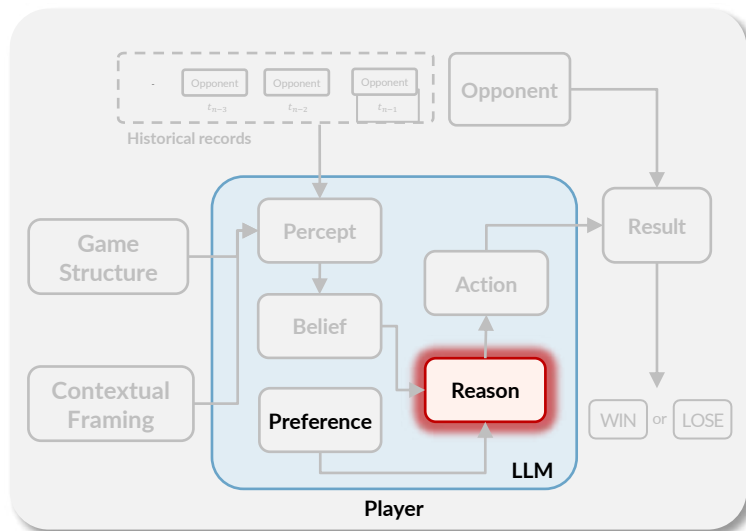
# Reasoning : Theory-of-Mind



## First-Order ToM Modelling

From my perspective, please infer several beliefs about the opponent's game pattern/preference for each round when holding different cards and the public card (if have).

## Second-Order ToM Modelling

From my perspective, please infer under what circumstances is the opponent likely to be influenced by my actions? Additionally, in what situations would the opponent make decisions based solely on their own hand? From the perspective of the opponent (he cannot observe my card but only action), please infer several beliefs about my game pattern/preference when holding different cards.

**The theory of mind (ToM) can enhance GPT's performance in imperfect information games.**

# Neural Theory-of-Mind


THEY DON'T KNOW THAT WE KNOW THEY KNOW WE KNOW.

**8 tasks and 31 abilities in social cognition**
**ToMBENCH: Benchmarking Theory of Mind in Large Language Models**
*Tsinghua University*

**Longer and clearer narrative**
**Explicit personality traits**
**OpenToM: A Comprehensive Benchmark for Evaluating Theory-of-Mind Reasoning Capabilities of Large Language Models**
*King's College London*
*Huawei London Research Centre*
*The Alan Turing Institute*

**Interactions**
Our results indicate that this capacity has **not yet emerged** in any manner.
**FANTOM: A Benchmark for Stress-testing Machine Theory of Mind in Interactions**
*Yejin Choi*

**Higher-Order ToM**
**HI-TOM: A Benchmark for Evaluating Higher-Order Theory of Mind Reasoning in Large Language Models**
*University of Michigan*
*Westlake University*

**Feb 23, 2024**          **Feb 14, 2024**          **Oct 31, 2023**          **Oct 25, 2023**

**Jan 1, 2023**
**Minding Language Models' (Lack of) Theory of Mind: A Plug-and-Play Multi-Character Belief Tracker**
*Yejin Choi*

**Apr 3, 2023**
**Neural Theory-of-Mind? On the Limits of Social Intelligence in Large LMs**
*Yejin Choi*

**"Static" Text**
- reporting bias
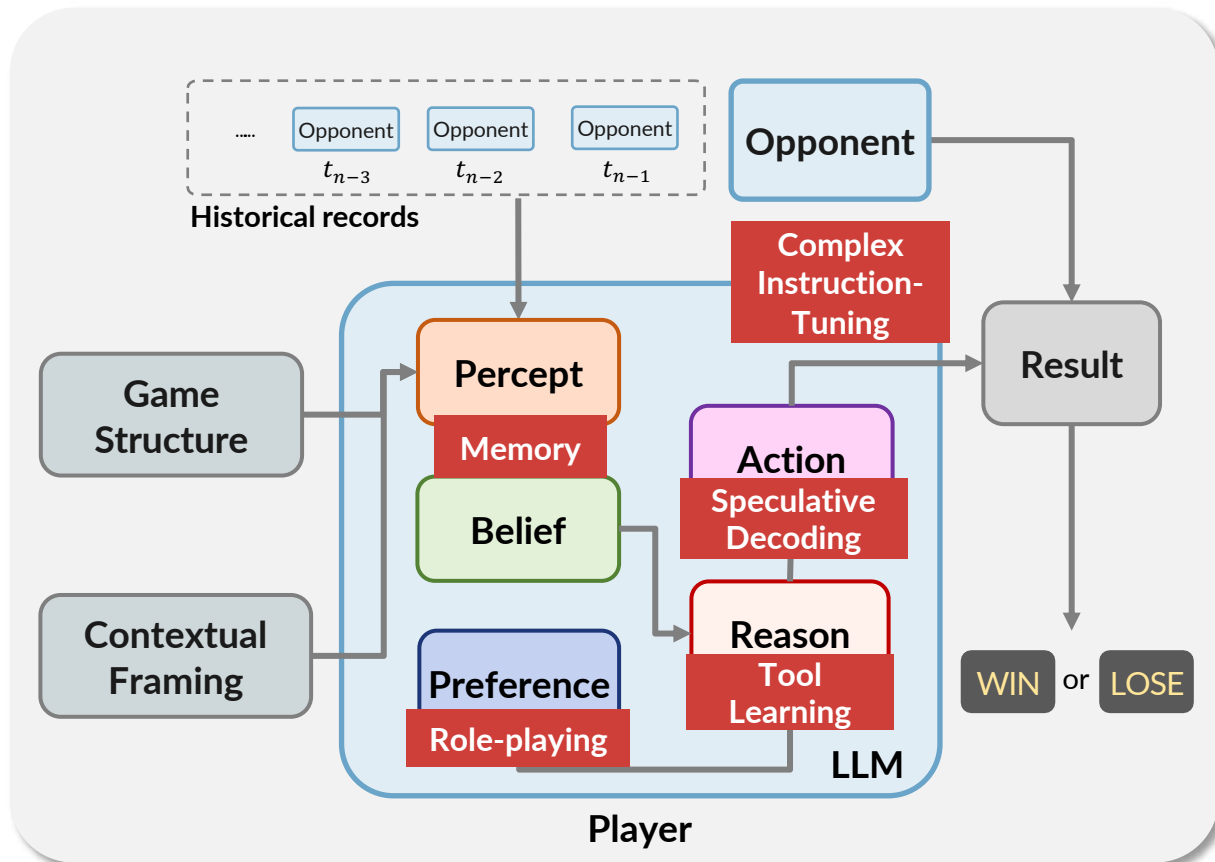- Lack of communicative intent and alternatives.
- Centering theory.

**May 24, 2023**
**Clever Hans or Neural Theory of Mind? Stress Testing Social Reasoning in Large Language Models**
*Yejin Choi*

**Oct 22, 2023**
**Theory of Mind for Multi-Agent Collaboration via Large Language Models**
*University of Pittsburgh*
*Carnegie Mellon University, Pittsburgh*
Evidence of **emergent** collaborative behaviors and high-order Theory of Mind capabilities among LLM-based agents.
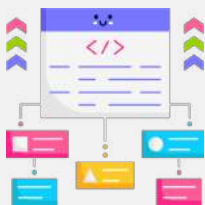
# Advanced Methods



Modified based on *Can Large Language Models Serve as Rational Players in Game Theory? A Systematic Analysis*

# Future Direction

**Unified Framework**

**Unified Metrics**

**Massive Experiments**

**Complex Scenarios**

**Capability Enhancement**

**Practical Applications**

# Takeaway

➢ Behavioral science for machines is of vital importance.
➢ Existing research utilizes game theory as a theoretical framework to investigate the strategic reasoning capabilities of large language models (LLMs).
➢ Preliminary experimental results indicate that while current LLMs possess some strategic reasoning abilities, these capabilities are not consistently stable.
➢ AI researchers and social science researchers need to communicate more frequently to enhance the depth of their research, including AI for Social Science and Social Science of AI.

# Thanks & QA

**Xiachong Feng**
**Postdoc Fellow**
**HKU**

**Haochuan Wang**
**Intern**
**HKU**

**Lingpeng Kong**
**Assistant Professor**
**HKU**

**Chuan Wu**
**Professor**
**HKU**